



LAWRENCE
LIVERMORE
NATIONAL
LABORATORY

LLNL-TR-419530

Dawn Usage, Scheduling, and Governance Model

S. Louis

November 5, 2009

Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

This work performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.



*ASC Sequoia
Acquisition Project*

Dawn Usage, Scheduling, and Governance Model

Version: 6.0 (FINAL VERSION)
Date: November 05, 2009
Document ID: LLNL-TR-419530



FINAL VERSION

Document Block

Item	Details	
Document Title	Dawn Usage, Scheduling, and Governance Model	
Document Type	Supporting Document for the Sequoia Acquisition Project	
Originator	Name:	Steve Louis
	Organization:	LLNL ICCD
	Contact Information:	stlouis@llnl.gov , 1-925-422-1550
Sponsor Point of Contact	Name:	Sander Lee
	Organization:	NNSA DP NA-114
	Contact Information:	Sander.Lee@nnsa.doe.gov , 1-202-586-2698
Document Tracking System	Fully qualified name:	LLNL-TR-419530
	Identifier	LLNL-TR-419530
	Version	Version 6.0
Last Revision	Version 6.0	
Submission Date	November 2, 2009	
Status	Final Version	
Release Date:	November 5, 2009	

Revision History

Revision Level	Date	Description	Change Summary
1.0	05/14/2009	Initial Draft	
2.0	05/27/2009	Revised Draft	Internal LLNL Modifications
3.0	06/17/2009	Revised Draft	Changes and Distribution to ASC HQ
4.0	09/17/2009	Revised Draft	Pools and Partitioning Changes
5.0	11/02/2009	Draft Final Version	OUO Designation Removed
6.0	11/05/2009	Final Version (R&R)	Review and Release Number Added

This work performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344.

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

FINAL VERSION

Dawn Usage, Scheduling, and Governance Model Advanced Simulation and Computing Program (Version 6 FINAL VERSION)

Steve Louis, LLNL

November 2, 2009



SUMMARY

This document describes Dawn use, scheduling, and governance concerns. Users started running full-machine science runs in early April 2009 during the initial open shakedown period. Scheduling Dawn while in the Open Computing Facility (OCF) was controlled and coordinated via phone calls, emails, and a small number of controlled banks. With Dawn moving to the Secure Computing Facility (SCF) in fall of 2009, a more detailed scheduling and governance model is required. The three major objectives are:

- Ensure Dawn resources are allocated on a program priority-driven basis
- Utilize Dawn resources on the job mixes for which they were intended
- Minimize idle cycles through use of partitions, banks and proper job mix

FINAL VERSION

The SCF workload for Dawn will be inherently different than Purple or BG/L, and therefore needs a different approach. Dawn's primary function is to permit adequate access for tri-lab code development in preparation for Sequoia, and in particular for weapons multi-physics codes in support of UQ. A second purpose is to provide time allocations for large-scale science runs and for UQ suite calculations to advance SSP program priorities.

This proposed governance model will be the basis for initial time allocation of Dawn computing resources for the science and UQ workloads that merit priority on this class of resource, either because they cannot be reasonably attempted on any other resources due to size of problem, or because of the unavailability of sizable allocations on other ASC capability or capacity platforms.

This proposed model intends to make the most effective use of Dawn as possible, but without being overly constrained by more formal proposal processes such as those now used for Purple CCCs.

ALLOCATION APPROACH

The proposed approach is to split the weekly schedule into two parts. The first half of the schedule would be from Monday 8:00am to Thursday 4:00pm (this is 80 hours per week, or 47.6% of total available hours) for code development and UQ calculations. The second half of the schedule would be from Thursday 4:00pm to Monday 8:00am (this is 88 hours per week, or 52.4% of total available hours) for large machine science runs and/or UQ.

During the Monday-to-Thursday "code development and UQ" portion of the week, allocations will be split evenly among the three labs for code development purposes and UQ workloads (28% of the time to each lab, with 16% reserved for systems and software development by LC). Because the job mix will vary significantly during this period, we propose to set up three different "pools" of Dawn nodes:

- 4 racks in a "pdebug" pool with 1 hour time limits (4,096 nodes total)
- 16 racks in a "pshort" pool with 4 hour time limits (16,384 nodes total)
- 16 racks in a "pbatch" pool with 12 hour time limits (16,384 nodes total)

During the Thursday-Monday "science run" portion of the week, allocations and "Dedicated Application Times" (DATs) will be determined primarily by programmatic priorities and decided, in advance, via email or telecons with Brian Carnes at LLNL as the point-of-contact, similar to how BG/L was scheduled in its early use phase. If no science runs or DATs are scheduled, a Monday-Thursday pdebug-pshort-pbatch pool scheme can remain in place. Because large science job mixes will tend to include mostly 32 or 36 rack jobs, we are proposing two pools of nodes during this half of the week:

- 4 racks in a "pdebug" pool with 1 hour time limits (4,096 nodes total)

FINAL VERSION

- 32 racks in a “pscience” pool with unlimited time limits (however, any science jobs still running would be terminated on Monday at 8:00am when scheduling reverts back to the Monday-Thursday pdebug-pshort-pbatch scenario)

The pscience pool is combined from the pshort and pbatch pools. If a full machine job needs all of Dawn, pscience and pdebug can be combined into a 36-rack “pfull” pool.

POOLS AND NODE PARTITIONING

pdebug

The pdebug pool would use a SLURM-only scheduling mechanism (no Moab resource management or bank prioritizations) that would effectively make the pdebug queue behave as a “first-in first-out” queue. The partitioning of the pdebug pool would allow overlapping partitioning schemes for 4x1K, 8*512, 16x256, and 32*128 node partitions. There would not be a 4K partition (which would tend to monopolize the entire pdebug pool, nor would there be any 2K partitions (which would tend to be sub-optimal due to mid-plane wiring considerations). Time limits for pdebug jobs would be 1 hour.

pshort

The pshort pool would use a SLURM + Moab fair-share scheduling mechanism, scheduling jobs based on bank structure and taking into consideration time used by recent jobs. The partitioning of the pshort pool would allow overlapping partitioning schemes for 16K, 2x8K, 4x4K, 8x2K, 16x1K, 32*512, 64x256, and 128*128 node partitions. Time limits for pshort jobs would be 4 hours to facilitate job throughput.

pbatch

The pbatch pool would use a SLURM + Moab fair-share scheduling mechanism, scheduling jobs based on bank structure and taking into consideration time used by recent jobs. The partitioning of the pbatch pool would allow overlapping partitioning schemes for 16K, 2x8K, 4x4K, 8x2K, 16x1K, 32*512, 64x256, and 128*128 node partitions. Time limits for pbatch jobs would be 12 hours to facilitate longer running jobs.

pscience and pfull

The pscience (32K node) and pfull (36K node) pools would also use a SLURM + Moab fair-share scheduling mechanism. However, projects wishing to utilize these pools at other than STANDBY would need allocations in the pscience or pfull banks. Controls over scheduling of science runs would be based on programmatic priorities and Dawn telecon or email decisions. The partitioning of pscience or pfull pools will be based on needs of the scheduled science or UQ run(s) to maximize Dawn node utilization.

FINAL VERSION

Note if a large science run needs less than the total nodes in the pscience or pfull pools, it is possible to run other jobs concurrently. These other jobs might be another science run, or might be multiple UQ runs. This allows the machine to accommodate both weapons science and UQ efforts simultaneously, so that both of the key future Sequoia missions can make progress simultaneously, without one having to wait for the other. This is the eventual intended mode of Sequoia, so it's be useful to test out this approach on Dawn.

A suitable bank structure will be developed to facilitate fair use of the machine. It should be noted that any time used over Thursday-Monday DAT periods will affect the fair share scheduling of jobs in pbatch and pshort during the Monday-Thursday period.

PROGRAMMATIC CONSIDERATIONS

We expect Dawn to be fairly shared among the three labs. DATs will be determined by priorities of the ASC Program and any milestones that may require Dawn computing resources. Schedules will be determined in view of milestones and priorities to ensure that Dawn use is tied strongly to documented program deliverables. Scheduling will take into account input from ASC Execs (or their designees), and input from other major science or UQ effort stakeholders as appropriate. For unanticipated or ad hoc urgent need, a means will be established for rapid changes to priorities or to the machine pool and partition configuration. DAT requests to use Dawn for large science runs (or large ensembles of UQ work) in the pscience or pfull pools during Thursday-Monday period should be prepared to define types of runs (i.e., size, duration, I/O, etc.) to be performed as well as the program deliverables expected as outcomes from those runs.

STANDBY JOB CONSIDERATIONS

Jobs submitted at STANDBY priority would be free to run if suitable partitions are available. STANDBY jobs will be immediately preempted (i.e., terminated) by any queued non-STANDBY calculation requiring nodes that are occupied by a STANDBY job. The STANDBY feature is now functional with SLURM/MOAB on Dawn.

ADMINISTRATIVE SUPPORT HOURS

We expect initial administrative support levels for Dawn when moved to the SCF to be the same as BG/L. This would mean during certain periods (after 8:00 p.m. and before 8:00am on weekdays, and after 5:00 p.m. and before 8:00am on weekends) there will not be any system administration support. This could result in delays for some operational procedures, including repartitioning to isolate bad nodes, or performing in-depth analysis for situations beyond simple hardware failure. Dawn users may be required to acclimate to lower levels of support off-hours, resulting in the possibility of idle partitions and a temporary inability to recover from failures. At some point (TBD) after Dawn General Availability (GA), support levels may, if necessary, increase to 24x7 similar to Purple.

SEQUOIA “SCALABLE APPLICATION PREPARATION” (SAP) PROJECT

FINAL VERSION

The SAP Project's mission is to support code teams in their preparation for Sequoia and will proactively assist tri-lab code developers to run on Dawn. The project will provide support and information, communicating with teams, LLNL staff, and IBM in such areas as compilers, debuggers, performance tools, simulators, third-party efforts, training, and documentation. The SAP team has already been instrumental in getting initial users onto Dawn for early science runs. The SAP Project Lead is Scott Futral, LLNL's Development Environment Group Lead and ASC ADEPT Project Co-lead (phone 1-925-422-1658).

DOCUMENTATION

A tutorial is available for new Dawn users. It begins with a brief history leading up to BG/P architecture. Dawn's configuration is presented, followed by information on the BG/P hardware architecture, including PowerPC 450, BG/P ASIC, double FPU, compute, I/O, login and service nodes, mid-planes, racks and the five BG/P networks. Topics relating to software development environments are covered, followed by usage information for BG/P compilers, MPI, OpenMP and Pthreads. Data alignment, math libraries, system configuration information, and specifics on running both batch and interactive jobs are presented. The tutorial concludes with a discussion about BG/P debugging and performance analysis tools.

Using the Dawn BG/P system:

<https://computing.llnl.gov/tutorials/bgp>

In addition to LLNL Dawn Tutorial web pages, other useful information can be found in IBM's published "Redbooks" on BG/P System Administration, BG/P Application Development, and Blue Gene Performance Analysis Tools. Redbooks of most probable interest to users are on application development and performance tools.

Blue Gene/P: System Administration:

<http://www.redbooks.ibm.com/abstracts/sg247417.html?Open>

Blue Gene/P: Application Development:

<http://www.redbooks.ibm.com/abstracts/sg247287.html?Open>

Blue Gene Solution: Performance Analysis Tools

<http://www.redbooks.ibm.com/abstracts/redp4256.html?Open>

These redbooks are available on the Dawn system in the /usr/local/docs directory, along with a text file named "dawn.basics" that has tips and other information for using Dawn.

DAWN SYSTEM ARCHITECTURE AND SYSTEM CONFIGURATION



Fig 1. LLNL Dawn BG/P Architecture

The picture above illustrates the basic architectural building blocks of Dawn. Each compute chip (one compute chip per compute node) contains four processing cores and threads. There is one compute chip (i.e., one node) on each compute card and I/O card, with 4 GB of memory per node. A compute node can be viewed as a compute chip plus its associated memory. There are 32 compute cards on each node card (making 32 nodes per node card). There are 16 node cards per mid-plane (a mid-plane is thus 512 nodes), and two mid-planes per cabinet (a cabinet has 1,024 nodes and 4,096 cores). The Dawn system has 36 cabinets for a total of 36,864 compute nodes and 147,456 PPC 450 cores.